# Structured Sparsity and Generalization

Andreas Maurer
Adalbertstr. 55
D-80799, München
*am@andreas-maurer.eu*

Massimiliano Pontil
Dept. of Computer Science, UCL
Gower St. London, UK
*m.pontil@cs.ucl.ac.uk*

September 5, 2011

### Abstract

We present a data dependent generalization bound for a large class of regularized algorithms which implement structured sparsity constraints. The bound can be applied to standard squared-norm regularization, the Lasso, the group Lasso, some versions of the group Lasso with overlapping groups, multiple kernel learning and other regularization schemes. In all these cases competitive results are obtained. A novel feature of our bound is that it can be applied in an infinite dimensional setting such as the Lasso in a separable Hilbert space or multiple kernel learning with a countable number of kernels.

## 1  Introduction

We study a class of regularization methods used to learn a linear function from a finite set of examples. The regularizer is expressed as an infimum convolution which involves a set $\mathcal{M}$ of linear transformations (see equation (1) below). As we shall see, this regularizer generalizes, depending on the choice of the set $\mathcal{M}$, the regularizers used by several learning algorithms, such as ridge regression, the Lasso, the group Lasso [22], multiple kernel learning [10], the group Lasso with overlap [6], and the regularizers in [16].

We give a bound on the Rademacher average of the linear function class associated with this regularizer. The result matches existing bounds in the above mentioned cases but also admits a novel, dimension free interpretation. In particular, the bound applies to the Lasso in $\ell_2$ or to multiple kernel learning with a countable number of kernels, under certain finite second-moment conditions.

Let $H$ be a real Hilbert space with inner product $\langle \cdot, \cdot \rangle$ and induced norm $\| \cdot \|$. Let $\mathcal{M}$ be a finite or countably infinite set of symmetric bounded linear operators on $H$ such that for every $x \in H$, $x \neq 0$, there is some linear operator $M \in \mathcal{M}$ with $Mx \neq 0$ and that $\sup_{M \in \mathcal{M}} \|\|M\|\| < \infty$, where $\|\| \cdot \|\|$ is the

operator norm. Define the function $\|\cdot\|_{\mathcal{M}} : H \to \mathbb{R}_+ \cup \{\infty\}$ by

$$\|\beta\|_{\mathcal{M}} = \inf \left\{ \sum_{M \in \mathcal{M}} \|v_M\| : v_M \in H, \ \sum_{M \in \mathcal{M}} M v_M = \beta \right\}. \tag{1}$$

It is shown in Section 3.2 that the chosen notation is justified, because $\|\cdot\|_{\mathcal{M}}$ is indeed a norm on the subspace of $H$ where it is finite, and the dual norm is, for every $z \in H$, given by

$$\|z\|_{\mathcal{M}*} = \sup_{M \in \mathcal{M}} \|M z\|.$$

The somewhat complicated definition of $\|\cdot\|_{\mathcal{M}}$ is contrasted by the simple form of the dual norm.

Using well known techniques, as described in [9] and [2], our study of generalization reduces to the search for a good bound on the empirical Rademacher complexity of a set of linear functionals with $\|\cdot\|_{\mathcal{M}}$-bounded weight vectors

$$\mathcal{R}_{\mathcal{M}}(\mathbf{x}) = \frac{2}{n} \mathbb{E} \sup_{\beta: \ \|\beta\|_{\mathcal{M}} \leq 1} \sum_{i=1}^{n} \epsilon_i \langle \beta, x_i \rangle, \tag{2}$$

where $\mathbf{x} = (x_1, \ldots, x_n) \in H^n$ is a sample vector representing observations, and $\epsilon_1, \ldots, \epsilon_n$ are Rademacher variables, mutually independent and each uniformly distributed on $\{-1, 1\}$[1]. Given a bound on $\mathcal{R}_{\mathcal{M}}(\mathbf{x})$ we obtain uniform bounds on the estimation error, for example using the following standard result (adapted from [2]), where the Lipschitz function $\phi$ is to be interpreted as a loss function.

**Theorem 1** *Let $\mathbf{X} = (X_1, \ldots, X_n)$ be a vector of iid random variables with values in $H$, let $X$ be iid to $X_1$, let $\phi : \mathbb{R} \to [0, 1]$ have Lipschitz constant $L$ and $\delta \in (0, 1)$. Then with probability at least $1 - \delta$ in the draw of $\mathbf{X}$ it holds, for every $\beta \in \mathbb{R}^d$ with $\|\beta\|_{\mathcal{M}} \leq 1$, that*

$$\mathbb{E}\phi\left(\langle \beta, X \rangle\right) \leq \frac{1}{n} \sum_{i=1}^{n} \phi\left(\langle \beta, X_i \rangle\right) + L \ \mathcal{R}_{\mathcal{M}}(\mathbf{X}) + \sqrt{\frac{9 \ln 2/\delta}{2n}}.$$

A similar (slightly better) bound is obtained if $\mathcal{R}_{\mathcal{M}}(\mathbf{X})$ is replaced by its expectation $\mathcal{R}_{\mathcal{M}} = \mathbb{E}\mathcal{R}_{\mathcal{M}}(\mathbf{X})$ (see [2]).

The following is the main result of this paper and leads to consistency proofs and finite sample generalization guarantees for all algorithms which use a regularizer of the form (1). A proof is given in Section 3.3.

---

[1]Our definition coincides with the one in [2], while other authors omit the factor of 2. This is relevant when comparing the constants in different bounds.

**Theorem 2** *Let* $\mathbf{x} = (x_1, \ldots, x_n) \in H^n$ *and* $\mathcal{R}_{\mathcal{M}}(\mathbf{x})$ *be defined as in (2). Then*

$$
\begin{aligned}
\mathcal{R}_{\mathcal{M}}(\mathbf{x}) \ \leq \ & \frac{2^{3/2}}{n} \sqrt{\sup_{M \in \mathcal{M}} \sum_{i=1}^{n} \|Mx_i\|^2} \left( 2 + \sqrt{\ln \left( \sum_{M \in \mathcal{M}} \frac{\sum_i \|Mx_i\|^2}{\sup_{N \in \mathcal{M}} \sum_j \|Nx_j\|^2} \right)} \right) \\
\leq \ & \frac{2^{3/2}}{n} \sqrt{\sum_{i=1}^{n} \|x_i\|_{\mathcal{M}*}^2} \left( 2 + \sqrt{\ln |\mathcal{M}|} \right).
\end{aligned}
$$

The second inequality follows from the first, the inequality

$$
\sup_{M \in \mathcal{M}} \sum_{i=1}^{n} \|Mx_i\|^2 \leq \sum_{i=1}^{n} \|x_i\|_{\mathcal{M}*}^2
$$

(a fact which will be tacitly used in the sequel) and the observation that every summand in the logarithm appearing in the first inequality is bounded by 1. Of course the second inequality is relevant only if $\mathcal{M}$ is finite. In this case we can draw the following conclusion: If we have an a priori bound on $\|X\|_{\mathcal{M}*}$ for some data distribution, say $\|X\|_{\mathcal{M}*} \leq C$, and $\mathbf{X} = (X_1, \ldots, X_n)$, with $X_i$ iid to $X$, then

$$
\mathcal{R}_{\mathcal{M}}(\mathbf{X}) \leq \frac{2^{3/2}C}{\sqrt{n}} \left( 2 + \sqrt{\ln |\mathcal{M}|} \right),
$$

thus passing from a data-dependent to a distribution dependent bound. In Section 2 we show that this recovers existing results for many regularization schemes.

But the first bound in Theorem 2 can be considerably smaller than the second and may be finite even if $\mathcal{M}$ is infinite. This gives rise to some appearantly novel features, even in the well studied case of the Lasso, when there is a (finite but potentially large) $\ell_2$-bound on the data.

**Corollary 3** *Under the conditions of Theorem 2 we have*

$$
\mathcal{R}_{\mathcal{M}}(\mathbf{x}) \leq \frac{2^{3/2}}{n} \sqrt{\sup_{M \in \mathcal{M}} \sum_i \|Mx_i\|^2} \left( 2 + \sqrt{\ln \frac{1}{n} \sum_i \sum_{M \in \mathcal{M}} \|Mx_i\|^2} \right) + \frac{2}{\sqrt{n}}.
$$

A proof is given in Section 3.3. To obtain a novel distribution dependent bound we retain the condition $\|X\|_{\mathcal{M}*} \leq C$ and replace finiteness of $\mathcal{M}$ by the condition that

$$
R^2 := \mathbb{E} \sum_{M \in \mathcal{M}} \|MX\|^2 < \infty. \tag{3}
$$

Taking the expectation in Corollary 3 then gives a bound on the expected Rademacher complexity

$$
\mathcal{R}_{\mathcal{M}} \leq \frac{2^{3/2}C}{\sqrt{n}} \left( 2 + \sqrt{\ln R^2} \right) + \frac{2}{\sqrt{n}}. \tag{4}
$$

The key features of this result are the dimension-independence and the only logarithmic dependence on $R^2$, which in many applications turns out to be simply $R^2 = \mathbb{E} \left\| X \right\|^2$.

The rest of the paper is organized as follows. In the next section we specialize our results to different regularizers. In Section 3 we present the proof of Theorem 2 as well as the proof of other results mentioned above. In Section 4 we discuss the extension of these results to the $\ell_q$ case. Finally we draw our conclusions and comment on future work.

## 2　Examples

Before giving the examples we mention a great simplification in the definition of the norm $\left\| \cdot \right\|_{\mathcal{M}}$ which occurs when the members of $\mathcal{M}$ have mutually orthogonal ranges. A simple argument, given in Proposition 9 below shows that in this case

$$\left\| \beta \right\|_{\mathcal{M}} = \sum_{M \in \mathcal{M}} \left\| M^+ \beta \right\|,$$

where $M^+$ is the pseudoinverse of $M$. If, *in addition*, every member of $\mathcal{M}$ is an orthogonal projection $P$, the norm further simplifies to

$$\left\| \beta \right\|_{\mathcal{M}} = \sum_{P \in \mathcal{M}} \left\| P \beta \right\|,$$

and the quantity $R^2$ occurring in the second moment condition (3) simplifies to

$$R^2 = \mathbb{E} \sum_{P \in \mathcal{M}} \left\| P X \right\|^2 = \mathbb{E} \left\| X \right\|^2.$$

For the remainder of this section $\mathbf{X} = (X_1, \ldots, X_n)$ will be a generic iid random vector of data points, $X_i \in H$, and $X$ will be a generic data variable, iid to $X_i$. If $H = \mathbb{R}^d$ we write $(X)_k$ for the $k$-th coordinate of $X$, not to be confused with $X_k$, which would be the $k$-th member of the vector $\mathbf{X}$.

### 2.1　The Euclidean Regularizer

In this simplest case we set $\mathcal{M} = \{I\}$, where $I$ is the identity operator on the Hilbert space $H$. Then $\left\| \beta \right\|_{\mathcal{M}} = \left\| \beta \right\|$, $\left\| z \right\|_{\mathcal{M}*} = \left\| z \right\|$, and the bound on the empirical Rademacher complexity becomes

$$\mathcal{R}_{\mathcal{M}} \left( \mathbf{x} \right) \leq \frac{2^{5/2}}{n} \sqrt{\sum_i \left\| x_i \right\|^2},$$

worse by a constant factor of $2^{3/2}$ than the corresponding result in [2], a tribute paid to the generality of our result.

4

## 2.2 The Lasso

Let us first assume that $H = \mathbb{R}^d$ is finite dimensional and set $\mathcal{M} = \{P_1, \ldots, P_d\}$ where $P_k$ is the orthogonal projection onto the 1-dimensional subspace generated by the basis vector $e_k$. All the above mentioned simplifications apply and we have $\|\beta\|_{\mathcal{M}} = \|\beta\|_1$ and $\|z\|_{\mathcal{M}*} = \|z\|_\infty$. Writing $x_{ik}$ for the $k$-th coordinate of a data-point $x_i$, the bound on $\mathcal{R}_{\mathcal{M}}(\mathbf{x})$ now reads

$$\mathcal{R}_{\mathcal{M}}(\mathbf{x}) \leq \frac{2^{3/2}}{n} \sqrt{\sum_i \|x_i\|_\infty^2} \left(2 + \sqrt{\ln d}\right).$$

If $\|X\|_\infty \leq 1$ almost surely we obtain

$$\mathcal{R}_{\mathcal{M}}(\mathbf{X}) \leq \frac{2^{3/2}}{\sqrt{n}} \left(2 + \sqrt{\ln d}\right),$$

which agrees with the bound in [8] on the dominant term (see also [2],[15]).

Our last bound is useless if $d \geq e^n$ or if $d$ is infinite. But whenever the norm of the data has finite second moments we can use Corollary 3 and (4) to obtain

$$\mathcal{R}_{\mathcal{M}}(\mathbf{X}) \leq \frac{2^{3/2}}{\sqrt{n}} \left(2 + \sqrt{\ln \mathbb{E} \|X\|_2^2}\right) + \frac{2}{\sqrt{n}}.$$

For nontrivial results $\mathbb{E} \|X\|^2$ only needs to be subexponential in $n$.

We remark that a similar condition to equation (3) for the Lasso, replacing the expectation with the supremum over $X$, has been considered within the context of elastic net regularization [4].

## 2.3 The Weighted Lasso

The Lasso assigns an equal penalty to all regression coefficients, while there may be a priori information on the respective significance of the different coordinates. For this reason different weightings have been proposed (see e.g. [20]). In our framework an appropriate set of operators is $\mathcal{M} = \{\alpha_1 P_1, \ldots, \alpha_k P_k, \ldots\}$, with $\alpha_k > 0$ where $\alpha_k^{-1}$ is the penalty weight associated with the $k$-th coordinate. Then

$$\|\beta\|_{\mathcal{M}} = \sum_k \alpha_k^{-1} |\beta_k|$$

and

$$\|z\|_{\mathcal{M}*} = \sup_k \alpha_k |z_k|.$$

To further illustrate the use of Corollary 3 let us assume that the underlying space $H$ is infinite dimensional (i.e. $H = \ell_2(\mathbb{N})$), and make the compensating assumption that $\alpha \in H$, i.e. $\sum_k \alpha_k^2 = R^2 < \infty$. For simplicity we also assume that $\sup_k \alpha_k \leq 1$. Then, if $\|X\|_\infty \leq 1$ almost surely, we have both $\|X\|_{\mathcal{M}*} \leq 1$ and $\sum_k \alpha_k^2 (X)_k^2 \leq R^2$. Again we obtain

$$\mathcal{R}_{\mathcal{M}}(\mathbf{X}) \leq \frac{2^{3/2}}{\sqrt{n}} \left(2 + \sqrt{\ln R^2}\right) + \frac{2}{\sqrt{n}}.$$

So in this case the second moment bound is enforced by the weighting sequence.

## 2.4 The Group Lasso

Let $H = \mathbb{R}^d$ and let $\{J_1, \ldots, J_r\}$ be a partition of the index set $\{1, \ldots, d\}$. We take $\mathcal{M} = \{P_{J_1}, \ldots, P_{J_r}\}$ where $P_{J_\ell} = \sum_{i \in J_\ell} P_i$ is the projection onto the subspace spanned by the basis vector $e_i$. The ranges of the $P_{J_\ell}$ then provide an orthogonal decomposition of $\mathbb{R}^d$ and the above mentioned simplifications also apply. We get

$$\|\beta\|_{\mathcal{M}} = \sum_{\ell=1}^r \|P_{J_\ell}\beta\|$$

and

$$\|z\|_{\mathcal{M}*} = \max_{\ell=1}^r \|P_{J_\ell}z\|.$$

The algorithm which uses $\|\beta\|_{\mathcal{M}}$ as a regularizer is called the group Lasso (see e.g. [22]). It encourages vectors $\beta$ whose support lies the union of a small number of groups $J_\ell$ of coordinate indices. If we know that $\|P_{J_\ell}X\| \leq 1$ almost surely for all $\ell \in \{1, \ldots, r\}$ then we get

$$\mathcal{R}_{\mathcal{M}}(\mathbf{X}) \leq \frac{2^{3/2}}{\sqrt{n}}\left(2 + \sqrt{\ln r}\right), \tag{5}$$

in complete symmetry with the Lasso and essentially the same as given in [8]. If $r$ is prohibitively large or if different penalties are desired for different groups, the same remarks apply as in the previous two sections. Just as in the case of the Lasso the second moment condition (3) translates to the simple form $\mathbb{E}\|X\|_2^2 < \infty$.

## 2.5 Overlapping Groups

In the previous examples the members of $\mathcal{M}$ always had mutually orthogonal ranges, which gave a simple appearance to the norm $\|\beta\|_{\mathcal{M}}$. If the ranges are not mutually orthogonal, the norm has a more complicated form. For example, in the group Lasso setting, if the groups $J_\ell$ cover $\{1, \ldots, d\}$, but are not disjoint, we obtain the regularizer of [6], given by

$$\Omega_{\text{overlap}}(\beta) = \inf\left\{\sum_{\ell=1}^r \|v_\ell\| : (v_\ell)_{jk} = 0 \text{ if } k \notin J_\ell \text{ and } \sum_{\ell=1}^r v_\ell = \beta\right\}.$$

If $\|P_{J_\ell}X_i\| \leq 1$ almost surely for all $\ell \in \{1, \ldots, r\}$ then the Rademacher complexity of the set of linear functionals with $\Omega_{\text{overlap}}(\beta) \leq 1$ is bounded as in (5), in complete equivalence to the bound for the group Lasso.

The same bound also holds for the class satisfying $\Omega_{\text{group}}(\beta) \leq 1$, where the function $\Omega_{\text{group}}$ is defined, for every $\beta \in \mathbb{R}^d$, as

$$\Omega_{\text{group}}(\beta) = \sum_{\ell=1}^r \|P_{J_\ell}\beta\|$$

which has been proposed by [7, 23]. To see this we only have to show that $\Omega_{\text{overlap}} \leq \Omega_{\text{group}}$ which is accomplished by generating a disjoint partition $\{J'_\ell\}_{\ell=1}^r$ where $J'_\ell \subseteq J_\ell$, writing $\beta = \sum_{\ell=1}^r P_{J'_\ell}\beta$ and realizing that $\left\|P_{J'_\ell}\beta\right\| \leq \|P_{J_\ell}\beta\|$. The bound obtained from this simple comparison may however be quite loose.

## 2.6 Regularizers Generated from Cones

Our next example considers structured sparsity regularizers as in [16]. Let $\Lambda$ be a nonempty subset of the open positive orthant in $\mathbb{R}^d$ and define a function $\Omega_\Lambda : \mathbb{R}^d \to \mathbb{R}$ by

$$\Omega_\Lambda(\beta) = \frac{1}{2} \inf_{\lambda \in \Lambda} \sum_{j=1}^d \left( \frac{\beta_j^2}{\lambda_j} + \lambda_j \right).$$

If $\Lambda$ is a convex cone, then it is shown in [17] that $\Omega_\Lambda$ is a norm and that the dual norm is given by

$$\|z\|_{\Lambda*} = \sup \left\{ \left( \sum_{j=1}^d \mu_j z_j^2 \right)^{1/2} : \mu_j = \lambda / \|\lambda\|_1 \text{ with } \lambda \in \Lambda \right\}.$$

The supremum in this formula is evidently attained on the set $\mathcal{E}(\Lambda)$ of extreme points of the closure of $\{\lambda / \|\lambda\|_1 : \lambda \in \Lambda\}$. For $\mu \in \mathcal{E}(\Lambda)$ let $M_\mu$ be the diagonal matrix with entries $\sqrt{\mu_j}\delta_{jl}$ and let $\mathcal{M}_\Lambda$ be the collection of matrices $\mathcal{M}_\Lambda = \{M_\mu : \mu \in \mathcal{E}(\Lambda)\}$. Then

$$\|z\|_{\Lambda*} = \sup_{M \in \mathcal{M}_\Lambda} \|Mz\|.$$

Clearly $\mathcal{M}_\Lambda$ is uniformly bounded in the operator norm, so if $\Lambda$ is a cone and $\mathcal{E}(\Lambda)$ is at most countable, then $\|\cdot\|_{\Lambda*} = \|\cdot\|_{\mathcal{M}*}$, $\Omega_\Lambda = \|\cdot\|_{\mathcal{M}*}$ and our bounds apply. If $\mathcal{E}(\Lambda)$ is finite and $\mathbf{x}$ is a sample then the Rademacher complexity of the class with $\Omega_\Lambda(\beta) \leq 1$ is bounded by

$$\frac{2^{3/2}}{n} \sqrt{\sum_i \|x_i\|_{\Lambda*}^2} \left( 2 + \sqrt{\ln |\mathcal{E}(\Lambda)|} \right).$$

## 2.7 Kernel Learning

This is the most general case to which the simplification applies: Suppose that $H$ is the direct sum $H = \oplus_{j \in \mathcal{J}} H_j$ of an at most countable number of Hilbert spaces $H_j$. We set $\mathcal{M} = \{P_j\}_{j \in \mathcal{J}}$, where $P_j : H \to H$ is the projection on $H_j$. Then

$$\|\beta\|_{\mathcal{M}} = \sum_{j \in \mathcal{J}} \|P_j\beta\|$$

and

$$\|z\|_{\mathcal{M}*} = \sup_{j \in \mathcal{J}} \|P_j z\|.$$

Such a situation arises in multiple kernel learning [10] or the nonparametric group Lasso [14] in the following way: One has an input space $\mathcal{X}$ and a collection $\{K_j\}_{j \in \mathcal{J}}$ of positive definite kernels $K_j : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$. Let $\phi_j : \mathcal{X} \to H_j$ be the feature map representation associated with kernel $K_j$, so that, for every $x, t \in \mathcal{X}$ $K_j(x, t) = \langle \phi_j(x), \phi_j(t) \rangle$ (for background on kernel methods see, for example, [19]).

Suppose that $\mathbf{x} = (x_1, \dots, x_n) \in \mathcal{X}^n$ is a sample. Define the kernel matrix $\mathbf{K}_j = (K_j(x_i, x_k))_{i,k=1}^n$. Using this notation the bound in Theorem 2 reads

$$\mathcal{R}\left((\phi(x_1), \dots, \phi(x_n))\right) \le \frac{2^{3/2}}{n} \sqrt{\sup_{j \in \mathcal{J}} \mathrm{tr}\mathbf{K}_j} \left(2 + \sqrt{\ln \frac{\sum_{j \in \mathcal{J}} \mathrm{tr}\mathbf{K}_j}{\sup_{j \in \mathcal{J}} \mathrm{tr}\mathbf{K}_j}}\right).$$

In particular, if $\mathcal{J}$ is finite and $K_j(x, x) \le 1$ for every $x \in \mathcal{X}$ and $j \in \mathcal{J}$, then the the bound reduces to

$$\frac{2^{3/2}}{\sqrt{n}} \left(2 + \sqrt{\ln |\mathcal{J}|}\right),$$

essentially in agreement with [3, 8, 21]. For infinite or prohibitively large $\mathcal{J}$ the second moment condition now becomes

$$\mathbb{E} \sum_{j \in \mathcal{J}} K_j(X, X) < \infty.$$

## 3 Proofs

We first give some notation and auxiliary results, then we prove the results announced in the introduction.

### 3.1 Notation and Auxiliary Results

The Hilbert space $H$ and the collection $M$ are fixed throughout the following, as is the sample size $n \in \mathbb{N}$.

Recall that $\|\cdot\|$ and $\langle \cdot, \cdot \rangle$ denote the norm and inner product in $H$, respectively. For a linear transformation $M : \mathbb{R}^n \to H$ the Hilbert-Schmidt norm is defined as

$$\|M\|_{HS} = \left(\sum_{i=1}^n \|Me_i\|^2\right)^{1/2}$$

where $\{e_i : i \in \mathbb{N}\}$ is the canonical basis of $\mathbb{R}^n$.

We use bold letters ($\mathbf{x}$, $\mathbf{X}$, $\boldsymbol{\epsilon}$, $\dots$) to denote $n$-tuples of objects, such as vectors or random variables.

Let $\mathcal{X}$ be any space. For $\mathbf{x} = (x_1, \ldots, x_n) \in \mathcal{X}^n$, $1 \le k \le n$ and $y \in \mathcal{X}$ we use $\mathbf{x}_{k \leftarrow y}$ to denote the object obtained from $\mathbf{x}$ by replacing the $k$-th coordinate of $\mathbf{x}$ with $y$. That is

$$\mathbf{x}_{k \leftarrow y} = (x_1, \ldots, x_{k-1}, y, x_{k+1}, \ldots, x_n).$$

The following concentration inequality, known as the bounded difference inequality (see McDiarmid [13]), goes back to the work of Hoeffding [5]. We only need it in the weak form stated below.

**Theorem 4** *Let $F : \mathcal{X}^n \to \mathbb{R}$ and write*

$$B^2 = \sum_{k=1}^{n} \sup_{y_1, y_2 \in \mathcal{X},\ \mathbf{x} \in \mathcal{X}^n} \left( F\left(\mathbf{x}_{k \leftarrow y_1}\right) - F\left(\mathbf{x}_{k \leftarrow y_2}\right) \right)^2.$$

*Let $\mathbf{X} = (X_1, \ldots, X_n)$ be a vector of independent random variables with values in $\mathcal{X}$, and let $\mathbf{X}'$ be iid to $\mathbf{X}$. Then for any $t > 0$*

$$\Pr\left\{ F\left(\mathbf{X}\right) > \mathbb{E} F\left(\mathbf{X}'\right) + t \right\} \le e^{-2t^2/B^2}.$$

Finally we need a simple lemma on the normal approximation:

**Lemma 5** *Let $a, \delta > 0$. Then*

$$\int_{\delta}^{\infty} \exp\left( \frac{-t^2}{2a^2} \right) dt \le \frac{a^2}{\delta} \exp\left( \frac{-\delta^2}{2a^2} \right).$$

**Proof.** For $t \ge \delta/a$ we have $1 \le at/\delta$. Thus

$$\int_{\delta}^{\infty} \exp\left( \frac{-t^2}{2a^2} \right) dt = a \int_{\delta/a}^{\infty} e^{-t^2/2} dt \le \frac{a^2}{\delta} \int_{\delta/a}^{\infty} t e^{-t^2/2} dt = \frac{a^2}{\delta} \exp\left( \frac{-\delta^2}{2a^2} \right).$$

∎

## 3.2   Properties of $\|\cdot\|_{\mathcal{M}}$ and Duality

We state again the general conditions on the set $\mathcal{M}$.

**Condition 6** *$\mathcal{M}$ is a finite or countably infinite set of symmetric bounded linear operators on a real separable Hilbert space $H$ such that:*

(a) *For every $x \in H$ with $x \ne 0$, there exists $M \in \mathcal{M}$ such that $Mx \ne 0$;*

(b) *$\sup_{M \in \mathcal{M}} |||M||| < \infty$, where $||| \cdot |||$ is the operator norm.*

Denote $\mathcal{V}(\mathcal{M}) = \{v : v = (v_M)_{M \in \mathcal{M}}, \ v_M \in H\}$, so the definition of $\|\beta\|_{\mathcal{M}}$ reads

$$\|\beta\|_{\mathcal{M}} = \inf \left\{ \sum_{M \in \mathcal{M}} \|v_M\| : v \in \mathcal{V}(\mathcal{M}) \text{ and } \sum_{M \in \mathcal{M}} M v_M = \beta \right\}.$$

**Theorem 7** *We have that:*

(i) $\|\cdot\|_{\mathcal{M}}$ *is positive homogeneous and subadditive on* $\ell_1(\mathcal{M})$;

(ii) $\ell_1(\mathcal{M})$ *is a dense subspace of* $H$. *If* $\mathcal{M}$ *is finite or* $H$ *is finite dimensional then* $\ell_1(\mathcal{M}) = H$;

(iii) $\|\cdot\|_{\mathcal{M}}$ *is a norm on* $\ell_1(\mathcal{M})$.

**Proof.** (i) Positive homogeneity of $\|\cdot\|_{\mathcal{M}}$ is clear. For subadditivity let $\beta, \gamma \in \ell_1(\mathcal{M})$. Let $\epsilon > 0$ be arbitrary and choose $w^\beta, w^\gamma \in \mathcal{V}(\mathcal{M})$ such that $\sum_{M \in \mathcal{M}} M w_M^\beta = \beta$, $\sum_{M \in \mathcal{M}} M w_M^\gamma = \gamma$, $\sum_{M \in \mathcal{M}} \left\| w_M^\beta \right\| \le \|\beta\|_{\mathcal{M}} + \epsilon$ and $\sum_{M \in \mathcal{M}} \|w_M^\gamma\| \le \|\gamma\|_{\mathcal{M}} + \epsilon$. Then $w^\beta + w^\gamma \in \mathcal{V}(\mathcal{M})$ and $\sum_{M \in \mathcal{M}} M \left( w^\beta + w^\gamma \right)_M = \beta + \gamma$. Thus $w^\beta + w^\gamma$ is in the feasable set for the definition of $\|\beta + \gamma\|_{\mathcal{M}}$ and

$$
\begin{aligned}
\|\beta + \gamma\|_{\mathcal{M}} &= \inf \left\{ \sum_{M \in \mathcal{M}} \|v_M\| : v \in \mathcal{V}(\mathcal{M}) \text{ and } \sum_{M \in \mathcal{M}} M v_M = \beta + \gamma \right\} \\
&\le \sum_{M \in \mathcal{M}} \|w_M^\beta + w_M^\gamma\| \\
&\le \sum_{M \in \mathcal{M}} \|w_M^\beta\| + \sum_{M \in \mathcal{M}} \|w_M^\gamma\| \le \|\beta\|_{\mathcal{M}} + \|\gamma\|_{\mathcal{M}} + 2\epsilon.
\end{aligned}
$$

Since $\epsilon$ was arbitrary subadditivity follows.

(ii) It follows from (i) that $\ell_1(\mathcal{M})$ is a linear subspace of $H$. Let $S$ be the set of finite linear combinations of the form

$$S = \left\{ \sum_{i=1}^{K} M_i v_i : K \in \mathbb{N}, \ M_i \in \mathcal{M}, \ v_i \in H \right\}.$$

Then $S$ is a linear subspace of $\ell_1(\mathcal{M})$ and contains all vectors of the form $MMv = M^2 v$ where $M \in \mathcal{M}$ and $v \in H$. If $x \in H$ is perpendicular to all of $S$ then for all $M \in \mathcal{M}$ we must have $x \perp MMx \iff Mx = 0$, which implies $x = 0$ by condition (a). This shows that $S$ and therefore also $\ell_1(\mathcal{M})$ are dense in $H$. The second assertion of (ii) is an easy consequence of the first.

(iii) Suppose $\beta \in \ell_1(\mathcal{M})$, $\beta \neq 0$ and $\beta = \sum_{\mathcal{M}} M v_M$ with $v \in \mathcal{V}(\mathcal{M})$.

$$0 \le \|\beta\| = \left\| \sum_{M \in \mathcal{M}} M v_M \right\| \le \sup_{M \in \mathcal{M}} |||M||| \sum_{M \in \mathcal{M}} \|v_M\|.$$

Taking the infimum on the right hand side we obtain that

$$\|\beta\|_{\mathcal{M}} \geq \frac{\|\beta\|}{\sup\limits_{M \in \mathcal{M}} \||M\||} > 0,$$

where condition (b) was used. Together with (i) this implies that $\|\cdot\|_{\mathcal{M}}$ is a norm on $\ell_1(\mathcal{M})$. ∎

From now on we refer to $\ell_1(\mathcal{M})$ as the normed linear space with norm $\|\cdot\|_{\mathcal{M}}$.

**Theorem 8** *Let $z \in H$. The linear functional $\beta \mapsto \langle \beta, z \rangle$ is bounded on $\ell_1(\mathcal{M})$ and has norm*

$$\|z\|_{\mathcal{M}*} = \sup\limits_{M \in \mathcal{M}} \|Mz\|.$$

**Proof.** Let $F$ be the dual norm. By definition

$$
\begin{aligned}
F(z) &= \inf\{s : s\,\|\beta\|_{\mathcal{M}} - \langle \beta, z \rangle \geq 0, \forall \beta \in H\} \\
&= \inf\left\{s : \sum_{M \in \mathcal{M}} (s\,\|v_M\| - \langle Mv_M, z \rangle) \geq 0,\ \forall v \in \mathcal{V}(\mathcal{M})\right\} \\
&= \inf\{s : s\,\|v\| - \langle Mv, z \rangle \geq 0,\ \forall v \in H, \forall M \in \mathcal{M}\} \\
&= \inf\{s : s \geq \langle v, Mz \rangle,\ \forall v \in H, \|v\| = 1, \forall M \in \mathcal{M}\} \\
&= \inf\{s : s \geq \|Mz\|,\ \forall M \in \mathcal{M}\} \\
&= \sup\limits_{M \in \mathcal{M}} \|Mz\| = \|z\|_{\mathcal{M}*}.
\end{aligned}
$$

∎

**Proposition 9** *If the ranges of the members of $\mathcal{M}$ are mutually orthogonal then for $\beta \in \ell_1(\mathcal{M})$*

$$\|\beta\|_{\mathcal{M}} = \sum_{M \in \mathcal{M}} \|M^+\beta\|,$$

*where $M^+$ is the pseudoinverse of $M$.*

**Proof.** The ranges of the members of $\mathcal{M}$ provide an orthogonal decomposition of $H$, so

$$\beta = \sum_{M \in \mathcal{M}} M\left(M^+\beta\right),$$

where we used the fact that $MM^+$ is the orthogonal projection onto the range of $M$. Taking $v_M = M^+\beta$ this implies that $\|\beta\|_{\mathcal{M}} \leq \sum_{M \in \mathcal{M}} \|M^+\beta\|$. On the other hand, if $\beta = \sum_{N \in \mathcal{M}} Nv_N$, then, applying $M^+$ to this identity we see that $M^+Mv_M = M^+\beta$ for all $M$, so

$$\sum_{M \in \mathcal{M}} \|v_M\| \geq \sum_{M \in \mathcal{M}} \|M^+Mv_M\| = \sum_{M \in \mathcal{M}} \|M^+\beta\|,$$

which shows the reverse inequality. ∎

## 3.3 Data and Distribution Dependent Bounds

We use the bounded difference inequality to derive a concentration inequality for linearly transformed random vectors.

**Lemma 10** *Let $\boldsymbol{\epsilon} = (\epsilon_1, \ldots, \epsilon_n)$ be a vector of independent real random variables with $-1 \le \epsilon_i \le 1$, and $\boldsymbol{\epsilon}'$ iid to $\boldsymbol{\epsilon}$. Suppose that $M$ is a linear transformation $M : \mathbb{R}^n \to H$.*

(i) *Then for $t > 0$ we have*

$$\Pr\left\{ \|M\boldsymbol{\epsilon}\| \ge \mathbb{E}\|M\boldsymbol{\epsilon}'\| + t \right\} \le \exp\left( \frac{-t^2}{2\|M\|_{HS}^2} \right).$$

(ii) *If $\boldsymbol{\epsilon}$ is orthonormal (satisfying $\mathbb{E}\epsilon_i\epsilon_j = \delta_{ij}$), then*

$$\mathbb{E}\|M\boldsymbol{\epsilon}\| \le \|M\|_{HS}. \tag{6}$$

*and, for every $r > 0$,*

$$\Pr\left\{ \|M\boldsymbol{\epsilon}\| > t \right\} \le e^{1/r} \exp\left( \frac{-t^2}{(2+r)\|M\|_{HS}^2} \right).$$

**Proof.** (i) Define $F : [-1, 1]^n \to \mathbb{R}$ by $F(\mathbf{x}) = \|M\mathbf{x}\|$. By the triangle inequality

$$\sum_{k=1}^{n} \sup_{y_1, y_2 \in [-1,1], \, \mathbf{x} \in [-1,1]^n} \left( F(\mathbf{x}_{k \leftarrow y_1}) - F(\mathbf{x}_{k \leftarrow y_2}) \right)^2$$

$$\le \sum_{k=1}^{n} \sup_{y_1, y_2 \in [-1,1], \, \mathbf{x} \in [-1,1]^n} \left\| M(\mathbf{x}_{k \leftarrow y_1} - \mathbf{x}_{k \leftarrow y_2}) \right\|^2$$

$$= \sum_{k=1}^{n} \sup_{y_1, y_2 \in [-1,1]} (y_1 - y_2)^2 \|Me_k\|^2$$

$$\le 4\|M\|_{HS}^2$$

The result now follows from the bounded difference inequality (Theorem 4).

(ii) If $\boldsymbol{\epsilon}$ is orthonormal then it follows from Jensen's inequality that

$$\mathbb{E}\|M\boldsymbol{\epsilon}\| \le \left( \mathbb{E}\left\| \sum_{i=1}^{n} \epsilon_i Me_i \right\|^2 \right)^{1/2} = \left( \sum_i \|Me_i\|^2 \right)^{1/2} = \|M\|_{HS}.$$

For the second assertion of (ii) first note that from calculus we get $(t-1)^2/2 - t^2/(2+r) \ge -1/r$ for all $t \in \mathbb{R}$. This implies that

$$e^{-(t-1)^2/2} \le e^{1/r} e^{-t^2/(2+r)}. \tag{7}$$

12

Since $1/r \geq 1/(2+r)$ the inequality to be proved is trivial for $t \leq \|M\|_{HS}$. If $t > \|M\|_{HS}$ then, using $\mathbb{E}\|M\boldsymbol{\epsilon}\| \leq \|M\|_{HS}$, we have $t - E\|M\boldsymbol{\epsilon}\| \geq t - \|M\|_{HS} > 0$, so by part (i) and (7) we obtain

$$
\begin{aligned}
\Pr\{\|M\boldsymbol{\epsilon}\| \geq t\} &= \Pr\{\|M\boldsymbol{\epsilon}\| \geq E\|M\boldsymbol{\epsilon}\| + (t - E\|M\boldsymbol{\epsilon}\|)\} \\
&\leq \exp\left(\frac{-(t - E\|M\boldsymbol{\epsilon}\|)^2}{2\|M\|_{HS}^2}\right) \leq \exp\left(\frac{-(t - \|M\|_{HS})^2}{2\|M\|_{HS}^2}\right) \\
&= \exp\left(\frac{-(t/\|M\|_{HS} - 1)^2}{2}\right) \leq e^{1/r} e^{-(t/\|M\|_{HS})^2/(2+r)} \\
&= e^{1/r} \exp\left(\frac{-t^2}{(2+r)\|M\|_{HS}^2}\right).
\end{aligned}
$$

∎

We now use integration by parts, a union bound and the above concentration inequality to derive a bound on the expectation of the supremum of the norms $\|M\boldsymbol{\epsilon}\|$. This is the essential step in the proof of Theorem 2. It is by no means a new technique, in fact it appears many times in the book by Ledoux and Talagrand [11], but compared to the combinatorial approach in [3] it seems more suited to the study of the problem at hand, and gives insights into the fine structure of the logarithmic factor appearing in bounds for Lasso-like methods.

**Lemma 11** *Let $\mathcal{M}$ be a finite or countably infinite set of linear transformations $M : \mathbb{R}^n \to H$ and $\boldsymbol{\epsilon} = (\epsilon_1, \ldots, \epsilon_n)$ a vector of orthonormal random variables (satisfying $\mathbb{E}\epsilon_i\epsilon_j = \delta_{ij}$) with values in $[-1, 1]$. Then*

$$
\mathbb{E} \sup_{M \in \mathcal{M}} \|M\boldsymbol{\epsilon}\| \leq \sqrt{2} \sup_{M \in \mathcal{M}} \|M\|_{HS} \left(2 + \sqrt{\ln \frac{\sum_{M \in \mathcal{M}} \|M\|_{HS}^2}{\sup_{M \in \mathcal{M}} \|M\|_{HS}^2}}\right).
$$

**Proof.** To lighten notation we abbreviate $\mathcal{M}_\infty := \sup_{M \in \mathcal{M}} \|M\|_{HS}$ below. We now use integration by parts

$$
\begin{aligned}
\mathbb{E} \sup_{M \in \mathcal{M}} \|M\boldsymbol{\epsilon}\| &= \int_0^\infty \Pr\left\{\sup_{M \in \mathcal{M}} \|M\boldsymbol{\epsilon}\| > t\right\} dt \\
&\leq \mathcal{M}_\infty + \delta + \int_{\mathcal{M}_\infty + \delta}^\infty \Pr\left\{\sup_{M \in \mathcal{M}} \|M\boldsymbol{\epsilon}\| > t\right\} dt \\
&\leq \mathcal{M}_\infty + \delta + \sum_{M \in \mathcal{M}} \int_{\mathcal{M}_\infty + \delta}^\infty \Pr\{\|M\boldsymbol{\epsilon}\| > t\} dt,
\end{aligned}
$$

where we have introduced a parameter $\delta \geq 0$. The first inequality above follows from the fact that probabilities never exceed 1, and the second from a union bound. Now for any $M \in \mathcal{M}$ we can make a change of variables and use (6),

13

which gives $\mathbb{E} \left\| M\boldsymbol{\epsilon} \right\| \leq \left\| M \right\|_{HS} \leq \mathcal{M}_{\infty}$, so that

$$
\begin{aligned}
\int_{\mathcal{M}_{\infty}+\delta}^{\infty} \Pr\left\{ \left\| M\boldsymbol{\epsilon} \right\| > t \right\} dt \quad &\leq \quad \int_{\delta}^{\infty} \Pr\left\{ \left\| M\boldsymbol{\epsilon} \right\| > \mathbb{E}\left\| M\boldsymbol{\epsilon} \right\| + t \right\} dt \\
&\leq \quad \int_{\delta}^{\infty} \exp\left( \frac{-t^2}{2\left\| M \right\|_{HS}^2} \right) dt \\
&\leq \quad \frac{\left\| M \right\|_{HS}^2}{\delta} \exp\left( \frac{-\delta^2}{2\left\| M \right\|_{HS}^2} \right),
\end{aligned}
$$

where the second inequality follows from Lemma 10 (i), and the third from Lemma 5. Substitution in the previous chain of inequalities and using Hoelder's inequality (in the $\ell_1/\ell_{\infty}$-version) give

$$
\mathbb{E} \sup_{M \in \mathcal{M}} \left\| M\boldsymbol{\epsilon} \right\| \leq \mathcal{M}_{\infty} + \delta + \frac{1}{\delta}\left( \sum_{M \in \mathcal{M}} \left\| M \right\|_{HS}^2 \right) \exp\left( \frac{-\delta^2}{2\mathcal{M}_{\infty}^2} \right). \tag{8}
$$

We now set

$$
\delta = \sqrt{ 2\ln\left( e\frac{\sum_{M \in \mathcal{M}} \left\| M \right\|_{HS}^2}{\mathcal{M}_{\infty}^2} \right) } \mathcal{M}_{\infty}.
$$

Then $\delta \geq 0$ as required. The substitution makes the last term in (8) smaller than $\mathcal{M}_{\infty} / \left( e\sqrt{2} \right)$, and since $1 + 1/\left( e\sqrt{2} \right) < \sqrt{2}$, we obtain

$$
\mathbb{E} \sup_{M \in \mathcal{M}} \left\| M\boldsymbol{\epsilon} \right\| \leq \sqrt{2}\mathcal{M}_{\infty}\left( 1 + \sqrt{ \ln\left( \frac{e\sum_{M \in \mathcal{M}} \left\| M \right\|_{HS}^2}{\mathcal{M}_{\infty}^2} \right) } \right).
$$

Finally we use $\sqrt{\ln es} \leq 1 + \sqrt{\ln s}$ for $s \geq 1$. $\blacksquare$

**Proof of Theorem 2.** Let $\boldsymbol{\epsilon} = \left( \epsilon_1, \ldots, \epsilon_n \right)$ be a vector of iid Rademacher variables. For $M \in \mathcal{M}$ we use $M\mathbf{x}$ to denote the linear transformation $M\mathbf{x} : \mathbb{R}^n \to H$ given by $\left( M\mathbf{x} \right)\mathbf{y} = \sum_i \left( Mx_i \right) y_i$. We have

$$
\mathcal{R}_{\mathcal{M}}\left( \mathbf{x} \right) = \frac{2}{n}\mathbb{E} \sup_{\beta: \left\| \beta \right\|_{\mathcal{M}} \leq 1} \left\langle \beta, \sum_{i=1}^{n} \epsilon_i x_i \right\rangle \leq \frac{2}{n}\mathbb{E} \left\| \sum_{i=1}^{n} \epsilon_i x_i \right\|_{\mathcal{M}*} = \frac{2}{n}\mathbb{E} \sup_{M \in \mathcal{M}} \left\| M\mathbf{x}\boldsymbol{\epsilon} \right\|.
$$

Applying Lemma 11 to the set of transformations $\mathcal{M}\mathbf{x} = \left\{ M\mathbf{x} : M \in \mathcal{M} \right\}$ gives

$$
\mathcal{R}_{\mathcal{M}}\left( \mathbf{x} \right) \leq \frac{2^{3/2}\sup_{M \in \mathcal{M}} \left\| M\mathbf{x} \right\|_{HS}}{n}\left( 2 + \sqrt{ \ln\frac{\sum_{M \in \mathcal{M}} \left\| M\mathbf{x} \right\|_{HS}^2}{\sup_{M \in \mathcal{M}} \left\| M\mathbf{x} \right\|_{HS}^2} } \right).
$$

Substitution of $\left\| M\mathbf{x} \right\|_{HS}^2 = \sum_i \left\| Mx_i \right\|^2$ gives the first inequality of Theorem 2 and

$$
\sup_{M \in \mathcal{M}} \left\| M\mathbf{x} \right\|_{HS}^2 \leq \sum_i \sup_{M \in \mathcal{M}} \left\| Mx_i \right\|^2 = \sum_i \left\| x_i \right\|_{*\mathcal{M}}^2
$$

14

gives the second inequality. ∎

**Proof of Corollary 3.** From calculus we find that $t \ln t \geq -1/e$ for all $t > 0$. For $A, B > 0$ and $n \in \mathbb{N}$ this implies that

$$A \ln \frac{B}{A} = n\left[(A/n)\ln(B/n) - (A/n)\ln(A/n)\right] \leq A \ln(B/n) + n/e. \quad (9)$$

Now multiply out the first inequality of Theorem 2 and use (9) with

$$A = \sup_{M \in \mathcal{M}} \sum_i \|Mx_i\|^2 \text{ and } B = \sum_{M \in \mathcal{M}} \sum_i \|Mx_i\|^2.$$

Finally use $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$ for $a, b > 0$ and the fact that $2^{3/2}/\sqrt{e} \leq 2$. ∎

# 4 Extension to the $\ell_q(\mathcal{M})$ Case

There is a rather obvious extension of our framework, which should be mentioned for completeness: Let $q$ and $p$ be conjugate exponents (i.e. $1/q + 1/p = 1$) and define

$$\|\beta\|_{\mathcal{M}_q} = \inf\left\{\left(\sum_{M \in \mathcal{M}} \|v_M\|^q\right)^{1/q} : v_M \in H \text{ and } \sum_{M \in \mathcal{M}} Mv_M = \beta\right\},$$

in analogy to (1). Then $\|\beta\|_{\mathcal{M}_q}$ is a norm and the dual norm is given by

$$\|z\|_{\mathcal{M}_{q*}} = \left(\sum_{M \in \mathcal{M}} \|Mz\|^p\right)^{1/p}.$$

The proof of these facts is omitted in this version of the paper. In the following we give a result, which can be applied to cases analogous to those in Section 2, where it recovers existing results up to constant multiplicative factors.

**Theorem 12** *Let $\mathbf{x}$ be a sample and $\mathcal{R}_{\mathcal{M}_q}(\mathbf{x})$ the empirical Rademacher complexity of the class of linear functions parameterized by $\beta$ with $\|\beta\|_{\mathcal{M}_q} \leq 1$. Then for $1 < q \leq 2$*

$$\mathcal{R}_{\mathcal{M}_q}(\mathbf{x}) \leq \frac{2^{3/2}}{n}\sqrt{\pi p \sum_i \|x_i\|_{\mathcal{M}_{q*}}^2}.$$

The proof is analogous to the proof of Theorem 2, but somewhat more straightforward.

**Lemma 13** *Let $\mathcal{M}$ be a finite or countably infinite set of linear transformations $M : \mathbb{R}^n \to H$ and $\boldsymbol{\epsilon} = (\epsilon_1, \ldots, \epsilon_n)$ a vector of orthonormal random variables (satisfying $\mathbb{E}\epsilon_i\epsilon_j = \delta_{ij}$) with values in $[-1, 1]$. Then for $p \geq 2$*

$$\mathbb{E}\left[\left(\sum_{M \in \mathcal{M}} \|M\boldsymbol{\epsilon}\|^p\right)^{1/p}\right] \leq \sqrt{2\pi p}\left(\sum_{M \in \mathcal{M}} \|M\|_{HS}^p\right)^{1/p}.$$

**Proof.** First note that by standard results on the absolute moments of the normal distribution

$$\int_0^\infty t^{p-1} \exp\left(\frac{-t^2}{2}\right) dt \leq \sqrt{\frac{\pi}{2}}(p-2)!! \leq \sqrt{\frac{\pi}{2}}(1 \cdot 3 \cdot \ldots \cdot p - 2) \leq \sqrt{\frac{\pi}{2}}p^{p/2-1},$$

so

$$\left(p\int_0^\infty t^{p-1} \exp\left(\frac{-t^2}{2}\right) dt\right)^{1/p} \leq \sqrt{\frac{\pi}{2}}^{1/p} p^{1/2} \leq \sqrt{\frac{\pi p}{2}}. \tag{10}$$

Jensen's inequality and integration by parts give

$$\mathbb{E}\left(\sum_{M \in \mathcal{M}} \|M\boldsymbol{\epsilon}\|^p\right)^{1/p} \leq \left(\sum_{M \in \mathcal{M}} \mathbb{E}\|M\boldsymbol{\epsilon}\|^p\right)^{1/p} = \left(\sum_{M \in \mathcal{M}} p\int_0^\infty \Pr\{\|M\boldsymbol{\epsilon}\| > t\} t^{p-1} dt\right)^{1/p}$$

$$\leq \left(2p\sum_{M \in \mathcal{M}} \int_0^\infty t^{p-1} \exp\left(\frac{-t^2}{4\|M\|_{HS}^2}\right) dt\right)^{1/p},$$

where Lemma 10 (ii) was used in the last step with $r = 2$. A change of variables $t \to t/\left(\sqrt{2}\|M\|_{HS}\right)$ gives

$$\mathbb{E}\left(\sum_{M \in \mathcal{M}} \|M\boldsymbol{\epsilon}\|^p\right)^{1/p} \leq \left(2p\int_0^\infty t^{p-1} \exp\left(\frac{-t^2}{2}\right) dt \sum_{M \in \mathcal{M}} 2^{p/2} \|M\|_{HS}^p\right)^{1/p}$$

$$\leq 2^{1/p+1/2}\sqrt{\frac{\pi p}{2}}\left(\sum_{M \in \mathcal{M}} \|M\|_{HS}^p\right)^{1/p},$$

where we use (10) in the last inequality. ∎

**Proof of Theorem 12.** As in the proof of Theorem 2 we proceed using duality and apply Lemma 13 to the set of transformations $\mathcal{M}\mathbf{x} = \{M\mathbf{x} : M \in \mathcal{M}\}$.

$$\mathcal{R}_{\mathcal{M}_q}(\mathbf{x}) \leq \frac{2}{n}\mathbb{E}\left\|\sum \epsilon_i x_i\right\|_{\mathcal{M}_{q*}} = \frac{2}{n}\mathbb{E}\left[\left(\sum_{M \in \mathcal{M}} \|M\mathbf{x}\boldsymbol{\epsilon}\|^p\right)^{1/p}\right]$$

$$\leq \frac{2}{n}\sqrt{2\pi p}\left(\sum_{M \in \mathcal{M}} \|M\mathbf{x}\|_{HS}^p\right)^{1/p} = \frac{2}{n}\sqrt{2\pi p\left(\sum_{M \in \mathcal{M}} \left(\sum_i \|Mx_i\|^2\right)^{p/2}\right)^{2/p}}$$

$$\leq \frac{2}{n}\sqrt{2\pi p\sum_i\left(\sum_{M \in \mathcal{M}} \left(\|Mx_i\|^2\right)^{p/2}\right)^{2/p}} = \frac{2^{3/2}}{n}\sqrt{\pi p\sum_i \|x_i\|_{\mathcal{M}_{p*}}^2},$$

where the last inequality is just the triangle inequality in $\ell_{p/2}$. ∎

# 5 Conclusion and Future Work

We have presented a bound on the Rademacher average for linear function classes described by infimum convolution norms which are associated with a class of bounded linear operators on a Hilbert space. We highlighted the generality of the approach and its dimension independent features.

When the bound is applied to specific cases ($\ell_2$, $\ell_1$, mixed $\ell_1/\ell_2$ norms) it recovers existing bounds (up to small changes in the constants). The bound is however more general and allows for the possibility to remove the "$\log d$" factor which appears in previous bounds. Specifically, we have shown that the bound can be applied in infinite dimensional settings, provided that the moment condition (3) is satisfied. We have also applied the bound to multiple kernel learning. While in the standard case the bound is only slightly worse in the constants, the bound is potentially smaller and applies to the more general case in which there is a countable set of kernels, provided the expectation of the sum of the kernels is bounded.

An interesting question is whether the bound presented is tight. As noted in [3] the "$\log d$" is unavoidable. This result immediately implies that our bound is also tight, since we may choose $R^2 = d$ in equation (3).

A potential future direction of research is the application of our results in the context of sparsity oracle inequalities. In particular, it would be interesting to modify the analysis in [12], in order to derive dimension independent bounds. Another interesting scenario is the combination of our analysis with metric entropy.

## Acknowledgements

# References

[1] F. R. Bach, G.R.G. Lanckriet and M. I. Jordan. Multiple kernels learning, conic duality, and the SMO algorithm. Proceedings of the Twenty-first International Conference on Machine Learning, 2004.

[2] P. L. Bartlett and S. Mendelson. Rademacher and Gaussian Complexities: Risk Bounds and Structural Results. *Journal of Machine Learning Research*, 3: 463–482, 2002.

[3] C. Cortes, M. Mohri, A. Rostamizadeh. Generalization bounds for learning kernels. In Proceedings of the Twenty-seventh International Conference on Machine Learning (ICML 2010), 2010.

[4] C. De Mol, E. De Vito, L. Rosasco. Elastic-net regularization in learning theory. *Journal of Complexity*, 25(2): 201–230, 2009.

[5] W. Hoeffding, Probability inequalities for sums of bounded random variables, *Journal of the American Statistical Association*, 58:13–30, 1963.

[6] L. Jacob, G. Obozinski, J. P. Vert. Group Lasso with overlap and graph Lasso. In Proceedings of the Twenty-sixth International Conference on Machine Learning (ICML 2009), pages 433–440, 2009.

[7] R. Jenatton, J.-Y. Audibert, F. Bach. Structured variable selection with sparsity inducing norms. *Arxiv preprint arXiv:0904.3523*, 2009.

[8] S. M. Kakade, S. Shalev-Shwartz, A. Tewari. Regularization techniques for learning with matrices. *ArXiv preprint arXiv0910.0610*, 2010.

[9] V. Koltchinskii and D. Panchenko, Empirical margin distributions and bounding the generalization error of combined classifiers, *The Annals of Statistics*, 30(1):1–50.

[10] G. R. G. Lanckriet and N. Cristianini and P. Bartlett and L. El Ghaoui and M. I. Jordan. Learning the kernel matrix with semi-definite programming. *Journal of Machine Learning Research*, 5:27–72, 2004.

[11] M. Ledoux, M. Talagrand. *Probability in Banach Spaces*, Springer, 1991.

[12] K. Lounici, M. Pontil, A.B. Tsybakov and S. van de Geer. Oracle inequalities and optimal inference under group sparsity. *Annals of Statistics* (to appear).

[13] C.McDiarmid. Concentration, in *Probabilistic Methods of Algorithmic Discrete Mathematics*", pages 195–248, Springer, 1998.

[14] L. Meier, S.A. van de Geer, and P. Bühlmann. High-dimensional additive modeling. *Annals of Statistics*, 37(6B):3779–3821, 2009.

[15] R. Meir and T. Zhang. Generalization error bounds for Bayesian mixture algorithms. *Journal of Machine Learning Research*, 4:839–860, 2003.

[16] C. A. Micchelli, J. M. Morales, M. Pontil. A family of penalty functions for structured sparsity. In *Advances in Neural Information Processing Systems 23*, J. Lafferty et al. Editors, 1612–1623, 2010.

[17] C. A. Micchelli, J. M. Morales, M. Pontil. Regularizers for structured sparsity. arXiv:1010.0556v2, 2011.

[18] C. A. Micchelli and M. Pontil. Feature space perspectives for learning the kernel. *Machine Learning*, 66:297–319, 2007.

[19] J. Shawe-Taylor and N. Cristianini. *Kernel Methods for Pattern Analysis*, Cambridge University Press, 2004.

[20] T. Shimamura, S. Imoto, R. Yamaguchi and S. Miyano. Weighted Lasso in graphical Gaussian modeling for large gene network estimation based on microarray data. *Genome Informatics*, 19:142–153, 2007.

[21] Y. Ying and C. Campbell. Generalization bounds for learning the kernel problem. In COLT, 2009.

[22] M. Yuan and Y. Lin. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society, Series B (Statistical Methodology)*, 68(1): 49–67, 2006.

[23] P. Zhao and G. Rocha and B. Yu. Grouped and hierarchical model selection through composite absolute penalties. *Annals of Statistics*, 37(6A): 3468–3497, 2009.